

AD-A083 903

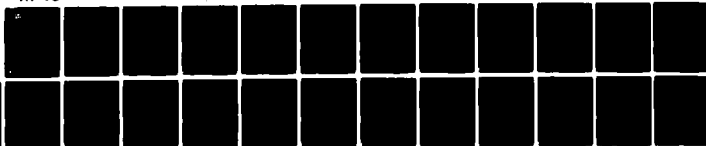
MISSOURI UNIV-COLUMBIA DEPT OF STATISTICS F/6 12/1  
MAXIMUM LIKELIHOOD ESTIMATION OF THE SURVIVAL FUNCTIONS OF STOC--ETC(U)  
MAR 80 R L DYKSTRA N00014-78-C-0655

UNCLASSIFIED

TR-91

NL

1 OF 1  
AD  
N00014-78-C-0655



END

DATE

FILED

DTIC



(12)  
R

**LEVEL**

AD A 0 8 3 9 0 3

University of Missouri-Columbia

**Maximum Likelihood Estimation of the  
Survival Functions of Stochastically  
Ordered Random Variables**

by

Richard L. Dykstra

Technical Report No. 91  
Department of Statistics

March 1980

Mathematical  
Sciences

DTIC  
ELECTE  
S MAY 7 1980

A

**DISTRIBUTION STATEMENT A**

Approved for public release  
Distribution Unlimited

80 5 5 018

DDC FILE COPY

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER 91	2. GOVT ACCESSION NO. AD-A083 903	3. RECIPIENT'S CATALOG NUMBER (9)
4. TITLE (and Subtitle) Maximum Likelihood Estimation of the Survival Functions of Stochastically Ordered Random Variables		5. TYPE OF REPORT & PERIOD COVERED Technical Report
7. AUTHOR(s) Richard L. Dykstra		6. PERFORMING ORG. REPORT NUMBER
9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Statistics University of Missouri Columbia, MO 65211		8. CONTRACT OR GRANT NUMBER(s) N00014-78-C-0655
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research Department of the Navy Arlington, VA		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) 12/26		12. REPORT DATE Mar 80
		13. NUMBER OF PAGES 21 pages
		15. SECURITY CLASS. (of this report) Unclassified
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNLIMITED.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) maximum likelihood equation; survival functions; stochastic ordering; censored observations; Kaplan-Meier product limit estimator; order restrictions		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) Many times populations exist which logically must satisfy a stochastic ordering requirement. Nevertheless, estimates of these populations may not bear out this stochastic ordering because of the inherent variability of the observations. This paper will consider the problem of finding maximum likelihood estimates of stochastically ordered survival functions		

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 68 IS OBSOLETE  
S N 0102-LF-014-6601

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

402600

HB

Unclassified  
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

20. (Abstract--continued)

for the cases a) one survival function being fixed in advance and b) estimating both survival functions when the data includes censored observations.

A numerical example is handled in detail to illustrate the solution to this problem.

AMS Classification numbers: Primary 60G05; Secondary 62N05

Maximum Likelihood Estimation of the Survival Functions  
of Stochastically Ordered Random Variables

by

Richard L. Dykstra

Technical Report No. 91  
March 1980

Prepared under contract N00014-78-C-0655  
for the Office of Naval Research

Reproduction in whole or in part is permitted for  
any purpose of the United States Government

Department of Statistics  
University of Missouri  
Columbia, Missouri 65211

Accession For	
NTIS G.A.I.	<input checked="checked" type="checkbox"/>
DDC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/	
Availability of	
Dist	Availability of
A	

MAXIMUM LIKELIHOOD ESTIMATION OF THE SURVIVAL FUNCTIONS  
OF STOCHASTICALLY ORDERED RANDOM VARIABLES

Richard L. Dykstra

ABSTRACT

Many times populations exist which logically must satisfy a stochastic ordering requirement. Nevertheless, estimates of these populations may not bear out this stochastic ordering because of the inherent variability of the observations. This paper will consider the problem of finding maximum likelihood estimates of stochastically ordered survival functions for the cases a) one survival function being fixed in advance and b) estimating both survival functions when the data includes censored observations.

A numerical example is handled in detail to illustrate the solution to this problem.

AMS Classification numbers: Primary 60G05; Secondary 62N05.

Key Words and Phrases: maximum likelihood estimation; survival functions; stochastic ordering; censored observations; Kaplan-Meier product limit estimator; order restrictions.

## I. INTRODUCTION

Many times a new population is created which can logically be only stochastically greater (or less) than the old population. Nevertheless, estimates of these populations may not bear out this stochastic ordering because of the inherent variability of the observations. Brunk, et.al. (1966) have given Maximum Likelihood Estimates (M.L.E.'s) of the C.D.F.'s of two stochastically ordered distributions when all observations are complete and the estimated distributions are required to be of the discrete type. This paper will consider the similar problem of finding M.L.E.'s of stochastically ordered survival functions for the cases a) one survival function being fixed in advance and b) estimating both survival functions when the data includes censored observations.

The unrestricted (except for the requirement of being a discrete distribution) M.L.E. of the survival function when censored data is involved has been developed by Kaplan and Meier (1958). Since the results of Kaplan and Meier are so well-known and widely used, this paper will largely conform to the notation developed there.

## II. NOTATION, THE PROBLEM, AND THE SOLUTION

Initially, we consider the following problem. Independent observations are taken from a discrete distribution on the positive part of the real line with survival function  $P(t)$ . We wish to find the M.L.E. of  $P(t)$  subject to  $P(t) \geq Q(t)$  for all  $t$  where  $Q$  denotes the survival function of a fixed discrete distribution with only a finite set of points possessing positive probability.

Complete observations (which we will call deaths) occur on a subset of the times  $S_1 < S_2 < \dots < S_m$  ( $S_0 = 0$  and  $S_{m+1} = \infty$  for convenience). The number of deaths at  $S_j$  is  $\delta_j$ . We let  $\lambda_j$  denote the number of censored observations (losses) in  $[S_j, S_{j+1})$ , assumed to occur at  $L_i^{(j)}$ ,  $i = 1, \dots, \lambda_j$ .

In Section III, we will prove that  $\hat{P}(t)$  may be expressed in the following manner. Let  $S_1 < S_2 < \dots < S_m$  denote the ordered values of the times of death combined with the points of positive probability under  $Q$ , and let  $n_i = \sum_{j=i}^m \delta_j + \lambda_j$  denote the number of items surviving just prior to  $S_i$ . Our  $p_j$  and  $q_j$  are related to  $P(\cdot)$  and  $Q(\cdot)$  respectively by

$$p_j = \ln [P(S_j)/P(S_{j-1})], \text{ and}$$

$$q_j = \ln [Q(S_j)/Q(S_{j-1})].$$

Then the restricted M.L.E. for  $t < S_m$  is given by

$$\hat{P}(t) = \exp \left[ \sum_{i; S_i \leq t} \hat{p}_i \right]$$

where  $\hat{p}_i$  is given in the following theorem.

(Throughout the paper we shall treat  $\infty$  and  $-\infty$  as real numbers greater than and less than all other real numbers respectively. We shall also adopt the conventions that



$\ln(0) = -\infty$  ,  $+\infty/+\infty = +1$  ,  $0(+\infty) = 0$  ,  $0/0 = 1$  , and  $0^0 = 1$  .)

Theorem 2.1. Let  $k_{a,b}^+$  denote the constant  $k$  such that

$$(2.1) \quad \sum_a^b \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right) = \sum_a^b q_j.$$

if  $k \geq 0$  and exists and 0 otherwise. Then if

$$(2.2) \quad \hat{k}_i = \min_{a \leq i} \max_{i \leq b} k_{a,b}^+,$$

$\hat{p}_i$  is expressible as

$$\hat{p}_i = \ln \left( \frac{n_i - \delta_i + \hat{k}_i}{n_i + \hat{k}_i} \right)$$

Of course if the last observation at say  $t^*$  corresponds to a loss, the  $\hat{p}(t)$  may be defined arbitrarily for  $t > t^*$  providing  $\hat{p}(t) \geq Q(t)$  and  $\hat{p}(t)$  is a survival function.

(Note that if  $Q(S_1) = Q(S_2) = \dots = Q(S_i) = 1$ , then

$$\sum_1^i \ln\left(\frac{n_j - S_j + k}{n_j + k}\right) = \sum_1^i q_j = 0$$

has the solution  $k = \infty$ . This leads to the intuitively reasonable solution  $\hat{p}_1 = \hat{p}_2 = \dots = \hat{p}_i = 0$ , even though the likelihood is identically zero in this case.)

Alternatively, the  $\hat{p}_j$ 's can be found more easily by the following algorithm:

1. Find the largest  $k_1 > 0$  such that

$$\sum_1^{i_1} \ln\left(\frac{n_j - \delta_j + k_1}{n_j + k_1}\right) = \sum_1^{i_1} q_j \text{ for } 0 < i_1 \leq m.$$

If more than one  $i_1$  works, choose largest.

2. Let  $\hat{p}_j = \ln\left(\frac{n_j - \delta_j + k_1}{n_j + k_1}\right)$  for  $1 \leq j \leq i_1$ .

3. Find the largest  $k_2 > 0$  such that

$$\sum_{i_1+1}^{i_2} \ln \left( \frac{n_j - \delta_j + k_2}{n_j + k_2} \right) = \sum_{i_1+1}^{i_2} q_j \text{ for } i_1 < i_2 \leq m.$$

If more than one  $i_2$  works, choose largest.

4. Let  $\hat{p}_j = \ln \left( \frac{n_j - \delta_j + k_2}{n_j + k_2} \right)$ ,  $i_1 < j < i_2$ .

5. Etc. If at some point, no such positive  $k_i$  exists, then

$$\hat{p}_j = \ln \left( \frac{n_j - \delta_j}{n_j} \right), \text{ for } i_{\ell-1} < j \leq m.$$

If the last observation at say  $t^*$  is a loss, then  $\hat{P}(t)$  may be may be defined arbitrarily beyond  $t^*$  as long as  $\hat{P}(t) \leq Q(t)$  and  $\hat{P}(t)$  is a survival function.

The order constraints  $P(t) \leq Q(t)$  can be handled similarly with the modifications given below.

If  $\delta_1 = \delta_2 = \dots = \delta_{i_0} = 0$  and  $\delta_{i_0+1} > 0$ , set  $\hat{p}_j = q_j$  for  $j = 1, 2, \dots, i_0$ . For  $j > i_0$ , Theorem 2.1 will still work by

- i) restricting  $a$  to be greater than  $i_0$ ,
- ii) defining  $k_{a,b}$  to be the value  $k$  in (2.1) if  $k < 0$  and exists, and 0 otherwise, and
- iii) interchanging the words min and max in (2.2).

The algorithm is modified to handle  $P(t) \leq Q(t)$  by

- i) setting  $\hat{p}_j = q_j$  for  $j = 1, 2, \dots, i_0$ ,
- ii) beginning the sum in step 1 at  $i_0+1$  rather than 1,
- iii) replacing  $k_i > 0$  by  $k_i < 0$ , and

- iv) replacing "largest  $k_i > 0$ " by "smallest  $k_i > 0$ " in steps 1 and 3 and replacing "positive" by "negative" in step 5.

The two-sample problem where a second set of independent observations are taken from a distribution with survival function  $Q(t)$  can be handled by essentially the same methods already given. In this case, let  $S_1 < \dots < S_m$  denote the ordered times of death of the combined samples. We let  $m_j$  denote the number of items in the second sample surviving just prior to  $S_j$ , and  $d_j$  the number of deaths at  $S_j$  in the second sample. Then Theorem 2.1 will still apply if (2.1) is replaced by

$$(2.1') \quad \sum_a^b \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right) = \sum_a^b \ln\left(\frac{m_j - d_j - k}{m_j - k}\right)$$

subject to a special case.

As is to be expected, when Theorem 2.1 applies

$$\hat{q}_j = \ln\left(\frac{m_j - d_j - \hat{k}_i}{m_j - \hat{k}_i}\right).$$

If  $d_1 = d_2 = \dots = d_{i_0} = 0$ ,  $d_{i_0+1} > 0$ , then we must set  $p_j = q_j$ ,  $j \leq i_0$ , in our likelihood function. It is straightforward to show that.

$$(2.3) \quad \hat{p}_j = \hat{q}_j = \frac{m_j + n_j - \delta_j}{m_j + n_j}, \quad j = 1, \dots, i_0$$

are the solutions for  $p_j$  and  $q_j$  in our equations. The earlier scheme then works if we require  $a$  to exceed  $i_0$ .

The algorithm will still work if we

i) define  $\hat{p}_j$  and  $\hat{q}_j$  as in (2.3) if  $j \leq i_0$

ii) replace

$$\sum_{i_{\ell-1}+1}^{i_\ell} q_j \text{ by } \sum_{i_{\ell-1}+1}^{i_\ell} \ln\left(\frac{m_j - d_j - k_\ell}{m_j - k_\ell}\right) \text{ for all } \ell$$

iii) and define

$$\hat{q}_j = \ln\left(\frac{m_j - d_j - k_l}{m_j - k_l}\right) \text{ for } i_{l-1} < j \leq i_l.$$

### III. DERIVATION FOR THE ONE SAMPLE PROBLEM

Clearly, the likelihood function of the observations is expressible as

$$(3.1) \quad L(P(t)) = \prod_{i=1}^{\lambda_0} P(L_i^{(0)}) \cdot \prod_{j=1}^m \{ [P(S_j - 0) - P(S_j)]^{\delta_j} \prod_{i=1}^{\lambda_j} P(L_i^{(j)}) \}$$

where  $S_1 < S_2 \dots < S_m$  are defined as in Section II. We wish to maximize this expression subject to the condition that  $P(t) \geq Q(t)$  for all  $t$ . Clearly to do this, the  $P(L_i^{(j)})$  and  $P(S_j - 0)$  should be as large as possible and the  $P(S_j)$  as small as possible. Decreasing the  $P(S_j)$  too much, however, may cause the order constraints to be violated. However, since we are dealing with discrete distribution, clearly it will suffice to maximize

$$\prod_{j=1}^m [P(S_{j-1}) - P(S_j)]^{\delta_j} P(S_j)^{\lambda_j}$$

subject to  $P(S_j) \geq Q(S_j)$ ,  $j = 0, \dots, m$ , since we can take  $P(\cdot)$  to be constant on the intervals  $[S_j, S_{j+1})$ .

Equivalently, we wish to maximize

$$\prod_{j=1}^m \left[ 1 - \frac{P(S_j)}{P(S_{j-1})} \right]^{\delta_j} \left[ \frac{P(S_{j-1})}{P(S_{j-2})} \dots P(S_1) \right]^{\delta_j} \left[ \frac{P(S_j)}{P(S_{j-1})} \dots P(S_1) \right]^{\lambda_j},$$

$$\text{or, letting } p'_j = \frac{P(S_j)}{P(S_{j-1})} \text{ and } q'_j = \frac{Q(S_j)}{Q(S_{j-1})},$$

to maximize

$$\prod_{j=1}^m [(1 - p'_j)^{\delta_j} p_j'^{\lambda_j} \cdot \prod_{i < j} p_i'^{\lambda_j + \delta_j}]$$

subject to

$$\prod_{j=1}^i p'_j \geq \prod_{j=1}^i q'_j \text{ for } i = 1, 2, \dots, m \text{ and } 1 \geq p'_j \geq 0 \text{ for all } j.$$

If we let  $n_i = \sum_{j=i}^m \delta_j + \lambda_j$  denote the number of items surviving just prior to  $S_i$ , our expression becomes

$$\prod_{j=1}^m (1 - p_j^{\delta_j}) p_j^{n_j - \delta_j}.$$

Finally, making the change of variables,

$$p_i = \ln p_i' \text{ and } q_i = \ln q_i'$$

and considering the natural log of the likelihood, our problem is to maximize

$$(3.2) \quad f(p_1, \dots, p_m) = \sum_{j=1}^m \delta_j \ln(1 - e^{p_j}) + (n_j - \delta_j) p_j$$

subject to the constraints

$$\sum_{j=1}^i p_j \leq \sum_{j=1}^i q_j \text{ and } 0 \geq p_i \geq -\infty \text{ for } i = 1, 2, \dots, m.$$

Suppose the constraint

$$\sum_{j=1}^i p_j = \sum_{j=1}^i q_j = c \quad (0 > c > -\infty)$$

is imposed when maximizing  $f(p_1, \dots, p_m)$ . Then by

writing  $p_i = c - \sum_{j=1}^{i-1} p_j$  and setting the partial derivatives equal to zero, we obtain the system of equations

$$(3.2) - \delta_j e^{p_j} / (1 - e^{p_j}) + n_j - \delta_j = 0, \quad j > i \text{ and}$$

$$(3.3) - \delta_j e^{p_j} / (1 - e^{p_j}) + n_j - \delta_j = -\delta_i e^{p_i} / (1 - e^{p_i}) + n_i - \delta_i = -k, \quad j < i.$$

Solving these equations gives the values

$$\hat{p}_j = \ln\left(\frac{n_j - \delta_j}{n_j}\right), \quad j > i, \text{ and}$$

$$\hat{p}_j = \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right), \quad j < i.$$

(Strictly speaking, equations (3.2) and (3.5) are not valid if  $\delta_j = 0$ . However, our solutions are still correct according to our previous conventions.) Moreover, noting that

$$\hat{p}_i = \ln\left(\frac{n_i - \delta_i + k}{n_i + k}\right),$$

by adding the first  $i$  equations it easily follows that  $k$  must be a real solution of the equation

$$\sum_{j=1}^i \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right) = \sum_{j=1}^i q_j = c.$$

In most situations,  $k$  will not have a closed form expression.

If  $m^*$  constraints are imposed in obtaining the M.L.E., say

$$\sum_{j=1}^{i_1} p_j = \sum_{j=1}^{i_1} q_j = c_1, \quad \dots, \quad \sum_{j=1}^{i_{m^*}} p_j = \sum_{j=1}^{i_{m^*}} q_j = c_{m^*},$$

then this is equivalent to the constraints

$$\sum_{j=1}^{i_1} p_j = c_1, \quad \sum_{j=i_1+1}^{i_2} p_j = c_2 - c_1, \quad \dots, \quad \sum_{j=i_{m^*}+1}^{i_{m^*}} p_j = c_{m^*} - c_{m^*-1}.$$

Each part of  $f(\cdot)$  may then be maximized separately to obtain a solution of the form

$$(3.4) \quad \hat{p}_j = \ln\left(\frac{n_j - \delta_j + k_\ell}{n_j + k_\ell}\right), \text{ if } i_{\ell-1} < j \leq i_\ell, \ell = 1, \dots, m^* + 1.$$

where  $k_\ell$  is the unique real solution of the equation

$$(3.5) \quad \sum_{i_{\ell-1}+1}^{i_\ell} \ln\left(\frac{n_j - \delta_j + k_\ell}{n_j + k_\ell}\right) = \sum_{i_{\ell-1}+1}^{i_\ell} q_j = c_\ell - c_{\ell-1}, \ell = 1, \dots, m^*.$$

(We take  $i_0, c_0$ , and  $k_{m^*+1}$  to all be zero.)

The problem of finding the stochastically ordered M.L.E. of the survival function is thus that of determining which constraints should be imposed.

To answer that question, the following lemmas are important.

Lemma 3.1. The function

$$f(p_1, p_2, \dots, p_m) = \sum_1^m \delta_j \ln(1 - e^{p_j}) + (n_j - \delta_j) p_j$$

is concave for  $0 \geq p_i \geq -\infty$ .

Proof. This easily follows since

$$\frac{\partial^2 f}{\partial p_i^2} = -\delta_i e^{p_i} / (1 - e^{p_i})^2 \leq 0.$$

Since  $f(p_1, \dots, p_m)$  is concave, and the constraints on the  $p_i$ 's are linear, the following lemma follows by arguments similar to those used in Theorem 1 of Dykstra and Madsen (1974).

Lemma 3.2. Assume that  $\Lambda$  and  $A$  denote column vectors,  $f(\cdot)$  is a concave real valued function defined on an appropriate subset of  $R^n$ , and that  $\Lambda_k$  maximizes  $f(\Lambda)$  subject to the constraints  $A_i' \Lambda \leq b_i$ ,  $i = 1, 2, \dots, k$ . Then if the additional constraint  $A_{k+1}' \Lambda \leq b_{k+1}$  is imposed, we may assume that

- a. if  $A_{k+1}' \Lambda_k \leq b_{k+1}$ ,  $\Lambda_{k+1} = \Lambda_k$ ;
- b. if  $A_{k+1}' \Lambda_k > b_{k+1}$ ,  $A_{k+1}' \Lambda_{k+1} = b_{k+1}$ .

The key in determining which constraints need be imposed in maximizing (3.2) is given in the following theorem.

Theorem 3.1. If the actual restricted M.L.E.'s are expressed in the form of (3.4)

then  $k_1 \geq k_2 \geq \dots \geq k_{m^*} \geq 0$ .

Proof. Without loss of generality, assume  $k_1 < k_2$ . Then if the maximum is found with the  $m^* - 1$  imposed constraints

$$\sum_{j=1}^{i_2} p_j = c_2, \dots, \sum_{j=1}^{i_{m^*}} p_j = c_{m^*},$$

the value of  $k_1^*$  which corresponds to the first constraint is such that  $k_1 < k_1^* < k_2$ . This easily follows since the  $\hat{p}_j$  are nondecreasing functions of the  $k_j$ . Since

$$\sum_{j=i}^{i_2} \ln\left(\frac{n_j - \delta_j + k_1^*}{n_j + k_1^*}\right) \leq \sum_{j=i}^{i_2} \ln\left(\frac{n_j - \delta_j + k_2}{n_j + k_2}\right)$$



for all  $i_1 < i \leq i_2$ , it follows that

$$\sum_{j=1}^i p_j^* \geq \sum_{j=1}^i q_j \quad \text{for } i = 1, \dots, m,$$

where the  $p_j^*$  denotes the values which maximize  $f(\cdot)$  imposing only the  $m^* - 1$  constraints. By Lemma 3.2, this implies that  $p_j^* = \hat{p}_j$  which is an obvious contradiction.

Theorem 3.2. The algorithm given in Section II obtains  $\hat{p}_j$ .

Proof. Straightforward from previous considerations.

The closed form expression for a particular  $\hat{p}_i$  given in Section II can also be obtained.

Proof (of Theorem 2.1). Clearly the expression is valid for  $\hat{p}_1$  by Theorem 3.2. Thus, using the notation of the algorithm given in Section II, it easily follows that  $k_1 = \hat{k}_1 \geq \hat{k}_2 \geq \dots \geq \hat{k}_m$ .

Let  $i$  denote the first integer  $\leq i_1$  such that  $\hat{k}_i = k_1$ . Then

$$\begin{aligned} \sum_{j=1}^{i_1} q_j &= \sum_{j=1}^{i_1} \ln\left(\frac{n_j - \delta_j + k_1}{n_j + k_1}\right) > \sum_{j=1}^{i-1} \ln\left(\frac{n_j - \delta_j + k_1}{n_j + k_1}\right) \\ &\quad + \sum_{j=i}^{i_1} \ln\left(\frac{n_j - \delta_j + \hat{k}_i}{n_j + \hat{k}_i}\right) \\ &\geq \sum_{j=1}^{i-1} q_j + \sum_{j=i}^{i_1} q_j \end{aligned}$$

which is a contradiction. Thus no such  $i$  exists and

$k_1 = \hat{k}_i$  for  $i = 1, 2, \dots, i_1$ . Essentially the same arguments handle subsequent intervals.

We should note that whenever  $\delta_j = 0$ , (no deaths are observed in the interval  $[S_j, S_{j+1})$ ),  $\hat{p}_j = 0$ , and the survival function  $\hat{P}(t)$  places zero probability over this interval.

Heuristically, the stochastically ordered M.L.E. acts like it has  $k_j$  additional items on test at the time  $S_j$ . Over the interval  $[S_j, S_{j+1})$ ,  $\hat{k}_j - \hat{k}_{j+1}$  of these items are lost, etc. Note, however, that the  $\hat{k}_j$  need not be integers.

#### IV. THE STOCHASTICALLY ORDERED TWO SAMPLE PROBLEM

Let us now consider the case where our observations, perhaps censored, come independently from two different discrete populations. As in Section III, we denote the true survival functions and the imposed ordering by

$$P(t) \geq Q(t) \text{ for all } t.$$

As before, we assume that  $P(0) = Q(0) = 1$  WLOG.

In a manner similar to that used in Section III, the function to be maximized may be expressed as

$$f(\underline{p}, \underline{q}) = \sum_{j=1}^m [\delta_j \ln(1 - e^{p_j}) + (n_j - \delta_j) p_j] + \\ [d_j \ln(1 - e^{q_j}) + (m_j - d_j) q_j]$$

where  $0 < S_1 < \dots < S_m$  denotes the ordered times of death of the combined sample,  $\delta_j$  ( $d_j$ ) denotes the number of deaths at  $S_j$  in the first (second) sample,  $n_j$  ( $m_j$ ) denotes the number of items surviving just prior to  $S_j$

in the first (second) sample, and  $p_j$  ( $q_j$ ) represents  $\ln[P(S_j)/P(S_{j-1})]$  ( $\ln[Q(S_j)/Q(S_{j-1})]$ ). The order constraints are still of the form

$$\sum_{j=1}^i p_j \geq \sum_{j=1}^i q_j \text{ for } i=1, \dots, m.$$

If we wish to maximize  $f$  subject to the one constraint

$$\sum_{j=1}^i p_j = \sum_{j=1}^i q_j,$$

then we may let  $q_i = \sum_{j=1}^i p_j - \sum_{j=1}^{i-1} q_j$ , and set the partial derivatives equal to zero. This results in the set of equations

$$\left. \begin{aligned} -\delta_j e^{p_j}/(1 - e^{p_j}) + n_j - \delta_j &= 0 \\ -d_j e^{q_j}/(1 - e^{q_j}) + m_j - d_j &= 0 \end{aligned} \right\} j > i$$

$$-\delta_j e^{p_j}/(1 - e^{p_j}) + n_j - \delta_j = d_i e^{q_i}/(1 - e^{q_i}) - m_i + d_i = k, j \leq i$$

$$-d_j e^{q_j}/(1 - e^{q_i}) + m_j - d_j = -d_i e^{q_i}/(1 - e^{q_i}) + m_i - d_i = j < i.$$

Solving these equations results in the solutions

$$\hat{p}_j = \ln\left(\frac{n_j - \delta_j}{n_j}\right), \quad j > i$$

$$\hat{q}_j = \ln\left(\frac{m_j - d_j}{m_j}\right), \quad j > i$$

$$\hat{p}_j = \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right), \quad j \leq i$$

$$\hat{q}_j = \ln\left(\frac{m_j - d_j - k}{m_j - k}\right), \quad j \leq i$$

where  $k$  is a real solution to the equation

$$\sum_{j=1}^i \ln\left(\frac{n_j - \delta_j + k}{n_j + k}\right) = \sum_{j=1}^i \ln\left(\frac{m_j - d_j - k}{m_j - k}\right).$$

If more constraints are imposed, similar solutions can be obtained for disjoint strings of the  $\hat{p}_j$  and  $\hat{q}_j$ .

Once again the key question is which constraints need be imposed in finding the true restricted M.L.E.'s. However, the same methods and lemmas used in proving Theorems 2.1, 3.1 and 3.2 will also suffice in the two-sample case and lead to the expressions given in Section II.

Another appealing method of handling the two-sample problem which will work if the samples do not include any censored data is to convert it to a one-sample problem as considered in Section 3 by making the observation that there must exist a survival function, say  $\hat{R}(t)$ , depending on the data such that

$$\hat{P}(t) \geq \hat{R}(t) \geq \hat{Q}(t) \text{ for all } t.$$

Intuitively, the more equality constraints imposed between  $\hat{P}(t)$  and  $\hat{Q}(t)$ , the more  $\hat{P}(t)$  will be pulled down and  $\hat{Q}(t)$  forced up (assuming the constraints  $\hat{P}(t) \geq \hat{Q}(t)$  are already imposed). The limiting case of this occurs if all equality constraints are imposed in which case

$$\hat{P}(t) = \hat{R}(t) = \hat{Q}(t)$$

where  $\hat{R}(t)$  denotes the Kaplan-Meier (1958) product limit estimator obtained from the pooled data. Thus if we construct  $\hat{R}(t)$ , treat it as fixed, and compute  $\hat{P}(t)$  as in Section III,  $\hat{P}(t)$  will be our restricted M.L.E. for  $\hat{P}(t)$ . With obvious changes,  $\hat{Q}(t)$  can be obtained similarly. Unfortunately, this method does not work when the samples include censored data although the end results are very close to the actual M.L.E.'s.

In the special case of no censored data explicit expressions for the  $k_{a,b}$  are given by

$$k_{a,b} = \frac{(\sum_a^b \delta_i) m_a - (\sum_a^b d_i) n_a}{\sum_a^b (d_i + \delta_i)}$$

This is of course a weighted average of  $m_a$  and  $-n_a$  with weight proportional to the number of deaths in the two populations for the appropriate interval. In this case the expressions in Section 2 are equivalent to those in Brunk, et al (1966).

## V. AN EXAMPLE

To illustrate the method, we apply it to some data gathered by Dr. Martin Alpert, Department of Cardiology of the University of Missouri Medical Center. Dr. Alpert's data consists of survival times for people who have had heart pacemakers implanted. We wish to estimate the survival functions separately for males and females, and will impose the constraint that the survival function for females never drops below that of males since it is well documented that females are longer lived. The data includes many censored observations of people who were lost to the study. The data on pages 18 and 19 has been coded to conform to the notation used in the paper.

The M.L.E.'s of ordered survival functions for males and females is shown on page 20.  $KM-P(T)$  and  $KM-Q(T)$  denote the unrestricted Kaplan-Meier estimates for the females and males respectively, while  $P(T)$  and  $Q(T)$  indicate the M.L.E.'s obtained using our algorithm. We note that  $P(T)$  has been forced up from  $KM-P(T)$  while  $Q(T)$  has been forced down from  $KM-Q(T)$ .

To try and get an idea of the overall effect of our order restrictions, we computed the expected values corresponding to our various survival functions when

they were all truncated at 97 months. These expectations were

$\frac{P(T)}{72.737}$	$\frac{Q(T)}{68.609}$	$\frac{KM-P(T)}{69.842}$	$\frac{KM-Q(T)}{70.837}$
-----------------------	-----------------------	--------------------------	--------------------------

Thus we see that our order restrictions increased the estimate of expected life (if truncated at 97 months) by nearly 3 months for females while decreasing it by approximately 2 months for males

A computer routine (in Fortran) for implementing this procedure is available upon request.

# PACEMAKER SURVIVAL DATA

Time (in months)	Males			Females		
	Losses in the interval $[S_j, S_{j+1})$	Deaths grouped to next higher month	No. of people surviving to just prior to $S_j$	Losses in the interval $(S_j, S_{j+1})$	Deaths grouped to next higher month	No. of people surviving to just prior to $S_j$
$S_j$	$\lambda_j$	$d_j$	$m_j$	$\lambda_j$	$\delta_j$	$n_j$
0	9	0	130	7	0	93
1	1	0	121	1	0	86
2	0	4	120	0	2	85
3	0	2	116	1	1	83
4	2	1	114	1	1	81
6	1	1	111	0	1	79
7	3	1	109	1	0	78
8	1	0	105	2	1	77
11	0	1	104	0	0	74
12	9	0	103	0	1	74
13	1	1	94	1	2	73
14	1	0	92	1	1	70
15	1	0	91	0	2	68
16	4	1	90	2	0	66
19	2	1	85	1	0	64
21	1	1	82	2	1	63
23	3	3	80	2	0	60
25	2	1	74	0	0	58
26	2	1	71	0	0	58
28	2	0	68	1	1	58

18



# PACEMAKER SURVIVAL DATA (Continued)

## Males

$S_j$	$\lambda_j$	$d_j$	$m_j$	$\lambda_j$	$\delta_j$	$n_j$
-------	-------------	-------	-------	-------------	------------	-------

30	2	0	66	3	1	56
34	0	1	64	1	0	52
35	2	0	63	1	1	51
37	1	1	61	1	0	49
38	1	0	59	0	1	48
39	7	1	58	1	1	47
43	0	1	50	2	0	45
45	2	0	49	1	1	43
47	0	0	47	1	1	41
48	2	1	47	2	0	39
51	2	1	44	0	1	37
54	2	0	41	0	1	36
55	1	1	39	0	0	35
56	0	1	37	2	0	35
57	1	1	36	1	0	33
58	6	0	34	3	1	32
62	0	1	28	0	0	28
63	1	1	27	2	0	28
65	2	0	25	0	1	26
67	3	0	23	3	1	25
71	1	0	20	2	1	21
73	1	1	19	1	0	18
76	0	1	17	1	0	17
78	1	1	16	0	0	16
81	1	0	14	2	1	16
86	2	0	13	3	1	13
89	6	0	11	1	1	9
97	4	1	5	7	0	7

TIME	P(T)	Q(T)	KM-P(T)	
1	0.100000000D 01	0.100000000D 01	0.100000000D 01	0.100000000D 01
2	0.901322987D 00	0.959148864D 00	0.976470580D 00	0.966666667D 00
3	0.971984481D 00	0.933723296D 00	0.964705882D 00	0.950000000D 00
4	0.962558383D 00	0.922510512D 00	0.952795933D 00	0.941666667D 00
6	0.953033005D 00	0.918068012D 00	0.940735225D 00	0.933133133D 00
7	0.953033005D 00	0.907505368D 00	0.940735225D 00	0.924621870D 00
8	0.943414522D 00	0.907505368D 00	0.928517885D 00	0.924621870D 00
11	0.943414522D 00	0.886426947D 00	0.928517885D 00	0.915731275D 00
12	0.933595828D 00	0.896426947D 00	0.915970346D 00	0.915731275D 00
13	0.913958440D 00	0.883962116D 00	0.898875260D 00	0.905909453D 00
14	0.904033118D 00	0.883962116D 00	0.878148478D 00	0.905909453D 00
15	0.883962116D 00	0.883962116D 00	0.858320580D 00	0.895829413D 00
16	0.883962116D 00	0.873750675D 00	0.858320580D 00	0.895829413D 00
19	0.883962116D 00	0.863458463D 00	0.858320580D 00	0.885332634D 00
21	0.870300550D 00	0.852704911D 00	0.838791604D 00	0.874835040D 00
23	0.870300550D 00	0.826032217D 00	0.838791604D 00	0.864178833D 00
25	0.870300550D 00	0.826032217D 00	0.838791604D 00	0.853741081D 00
26	0.870300550D 00	0.797012174D 00	0.838791604D 00	0.843716887D 00
28	0.855723933D 00	0.797012174D 00	0.826329758D 00	0.843716887D 00
30	0.840894599D 00	0.797012174D 00	0.808560584D 00	0.843716887D 00
34	0.840894599D 00	0.784224940D 00	0.808560584D 00	0.835924412D 00
35	0.824939002D 00	0.784224940D 00	0.793734886D 00	0.835924412D 00
37	0.824939002D 00	0.770999162D 00	0.793734886D 00	0.792712536D 00
38	0.808343023D 00	0.770999162D 00	0.777198743D 00	0.792712536D 00
39	0.791746244D 00	0.757303513D 00	0.760662599D 00	0.779045079D 00
43	0.791746244D 00	0.741622788D 00	0.760662599D 00	0.763464177D 00
45	0.774035703D 00	0.741622788D 00	0.742972771D 00	0.763464177D 00
47	0.755910451D 00	0.741622788D 00	0.724851484D 00	0.763464177D 00
48	0.755910451D 00	0.725249634D 00	0.724851484D 00	0.747220259D 00
51	0.736380310D 00	0.708102350D 00	0.705260903D 00	0.730237930D 00
54	0.716850170D 00	0.708102350D 00	0.685670323D 00	0.730237930D 00
55	0.716850170D 00	0.689115923D 00	0.685670323D 00	0.711513023D 00
56	0.716850170D 00	0.669591563D 00	0.685670323D 00	0.692203303D 00
57	0.716850170D 00	0.650067203D 00	0.685670323D 00	0.673053717D 00
58	0.695581689D 00	0.650067203D 00	0.664243125D 00	0.673053717D 00
62	0.695581689D 00	0.625315294D 00	0.664243125D 00	0.644901608D 00
63	0.695581689D 00	0.600623385D 00	0.664243125D 00	0.624978451D 00
65	0.670474300D 00	0.600623385D 00	0.632695313D 00	0.624978451D 00
67	0.645367711D 00	0.600623385D 00	0.613147500D 00	0.624978451D 00
71	0.616943634D 00	0.600623385D 00	0.583250000D 00	0.624978451D 00
73	0.616943634D 00	0.555895607D 00	0.583950000D 00	0.592084849D 00
76	0.616943634D 00	0.528897326D 00	0.583950000D 00	0.557256328D 00
78	0.616943634D 00	0.491899045D 00	0.583950000D 00	0.522427808D 00
81	0.582097538D 00	0.491899045D 00	0.567453125D 00	0.522427808D 00
86	0.542312029D 00	0.491899045D 00	0.505341346D 00	0.522427808D 00
89	0.491832769D 00	0.491899045D 00	0.449192308D 00	0.522427808D 00
97	0.491832769D 00	0.393519236D 00	0.449192308D 00	0.417942246D 00

## REFERENCES

1. Brunk, H. D., Franck, W. E., Hanson, D. L., and Hogg, R. V. (1966). Maximum likelihood estimation of the distributions of two stochastically ordered random variables. Journal of Amer. Statistical Assn. 61, 1067-1080.
2. Dykstra, R. L. and Madsen, R. W. (1976). Restricted maximum likelihood estimators for Poisson Parameters. Journal of Amer. Statistical Assn. 71, 711-718.
3. Kaplan, E. L. and Meier, Paul (1958). Nonparametric estimation from incomplete observations. Journal of Amer. Statistical Assn. 53, 457-481.